# Decision Trees

David Suffolk

5/3/2020

## Decision Trees

**Import Data**

```
library(rpart)
library(rpart.plot)
library(ggplot2)
data <- read.csv("./Project 02/nwCrow_bloodParasites_alaska_smith_2007_2008/nwCrow_sampling_alaska_smit
data_2 <- read.csv("./Project 02/nwCrow_bloodParasites_alaska_smith_2007_2008/nwCrow_bloodParasites_alas
de <- merge(data, data_2, by=0, all=TRUE)
head(de)
```

```
##   Row.names Field.ID      DATE  LOC   LAT    LONG SEX AGE AKD TARSUS WING MASS
## 1         1    75001 3/20/2007 SEWA 60.11 -149.44   1   1   1   55.8  283  448
## 2        10    75010 3/22/2007 KENA 60.55 -151.23   2   1   0   48.6  271  390
## 3       100    75100 3/12/2008 VALD 61.12 -146.35   2   1   0   44.6  264  317
## 4       101    86701 3/12/2008 VALD 61.12 -146.35   2   1   0   47.1  269  343
## 5       102    86702 3/12/2008 VALD 61.12 -146.35   2   2   0   52.2  291  415
## 6       103    86703 3/12/2008 VALD 61.12 -146.35   1   2   0   47.0  266  325
##   Extraction.. LEUC1 LEUC2 HAEM1 HAEM2 PLAS1 PLAS2 Leuc_GenBank_Accession
## 1      NOCR001     0     0     0     0     0     0
## 2      NOCR010     0     0     0     0     0     0
## 3      NOCR100     1     1     0     0     0     0                MG765394
## 4      NOCR101     0     0     0     0     0     0
## 5      NOCR102     1     0     0     0     0     0                MG765394
## 6      NOCR103     1     1     0     0     0     0                MG765394
##   Haem_GenBank_Accession Plas_GenBank_Accession
## 1
## 2
## 3
## 4
## 5
## 6
```

**One Hot Encoding**

```r
for(unique_value in unique(de$LOC)){

de[paste("LOC", unique_value, sep = ".")] <- ifelse(de$LOC == unique_value, 1, 0)

}
head(de)
```

```
##   Row.names Field.ID      DATE  LOC   LAT    LONG SEX AGE AKD TARSUS WING MASS
## 1         1    75001 3/20/2007 SEWA 60.11 -149.44   1   1   1   55.8  283  448
## 2        10    75010 3/22/2007 KENA 60.55 -151.23   2   1   0   48.6  271  390
## 3       100    75100 3/12/2008 VALD 61.12 -146.35   2   1   0   44.6  264  317
## 4       101    86701 3/12/2008 VALD 61.12 -146.35   2   1   0   47.1  269  343
## 5       102    86702 3/12/2008 VALD 61.12 -146.35   2   2   0   52.2  291  415
## 6       103    86703 3/12/2008 VALD 61.12 -146.35   1   2   0   47.0  266  325
##   Extraction.. LEUC1 LEUC2 HAEM1 HAEM2 PLAS1 PLAS2 Leuc_GenBank_Accession
## 1      NOCR001     0     0     0     0     0     0
## 2      NOCR010     0     0     0     0     0     0
## 3      NOCR100     1     1     0     0     0     0                MG765394
## 4      NOCR101     0     0     0     0     0     0
## 5      NOCR102     1     0     0     0     0     0                MG765394
## 6      NOCR103     1     1     0     0     0     0                MG765394
##   Haem_GenBank_Accession Plas_GenBank_Accession LOC.SEWA LOC.KENA LOC.VALD
## 1                                                      1        0        0
## 2                                                      0        1        0
## 3                                                      0        0        1
## 4                                                      0        0        1
## 5                                                      0        0        1
## 6                                                      0        0        1
##   LOC.HAIN LOC.JUNE LOC.HOME
## 1        0        0        0
## 2        0        0        0
## 3        0        0        0
## 4        0        0        0
## 5        0        0        0
## 6        0        0        0
```

**Filter Columns**

```r
de <- de[,c(7,8,9,10,11,12,14,16,18,23,24,25,26,27,28)]
head(de)
```

```
##   SEX AGE AKD TARSUS WING MASS LEUC1 HAEM1 PLAS1 LOC.SEWA LOC.KENA LOC.VALD
## 1   1   1   1   55.8  283  448     0     0     0        1        0        0
## 2   2   1   0   48.6  271  390     0     0     0        0        1        0
## 3   2   1   0   44.6  264  317     1     0     0        0        0        1
## 4   2   1   0   47.1  269  343     0     0     0        0        0        1
## 5   2   2   0   52.2  291  415     1     0     0        0        0        1
## 6   1   2   0   47.0  266  325     1     0     0        0        0        1
##   LOC.HAIN LOC.JUNE LOC.HOME
## 1        0        0        0
```
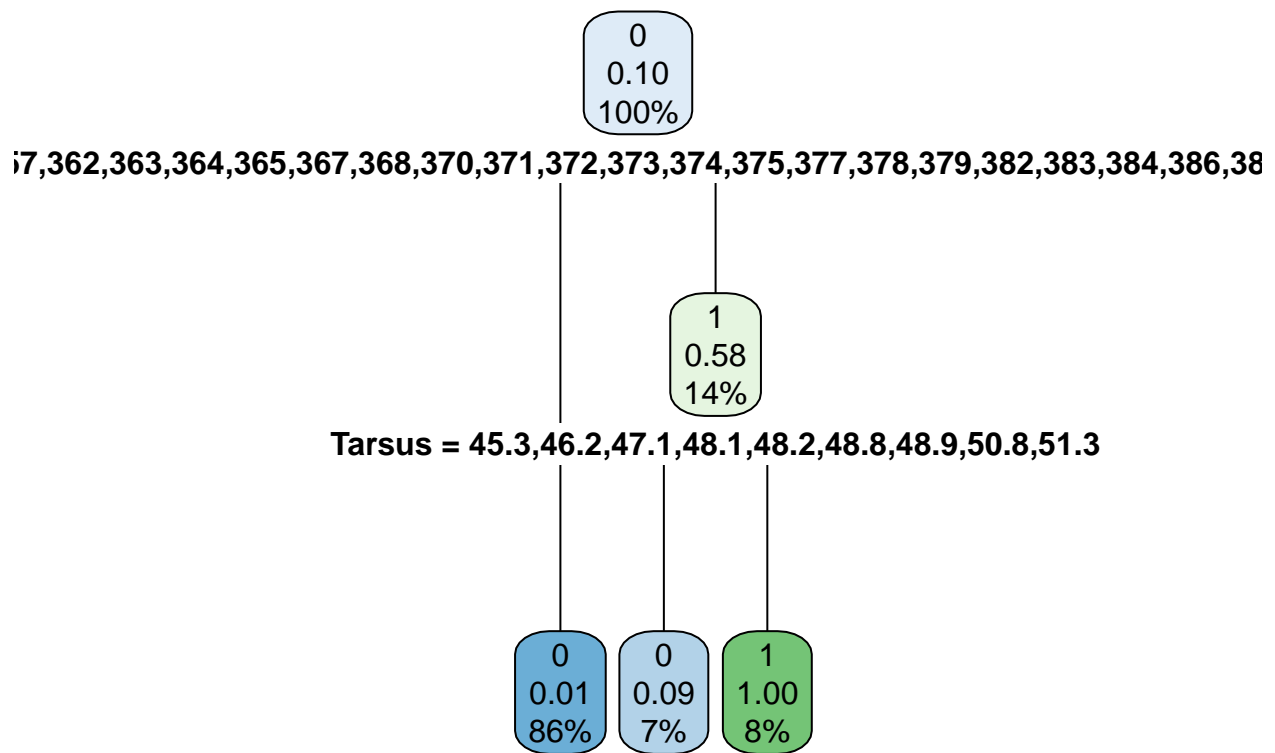
```
## 2          0          0          0
## 3          0          0          0
## 4          0          0          0
## 5          0          0          0
## 6          0          0          0
```

**Factoring**

```r
de$SEX <- as.factor(de$SEX)
de$AGE <- as.factor(de$AGE)
de$AKD <- as.factor(de$AKD)
de$TARSUS <- as.factor(de$TARSUS)
de$WING <- as.factor(de$WING)
de$MASS <- as.factor(de$MASS)
de$LEUC1 <- as.factor(de$LEUC1)
de$HAEM1 <- as.factor(de$HAEM1)
de$PLAS1 <- as.factor(de$PLAS1)
de$LOC.SEWA <- as.factor(de$LOC.SEWA)
de$LOC.KENA <- as.factor(de$LOC.KENA)
de$LOC.VALD <- as.factor(de$LOC.VALD)
de$LOC.HAIN <- as.factor(de$LOC.HAIN)
de$LOC.JUNE <- as.factor(de$LOC.JUNE)
de$LOC.HAIN <- as.factor(de$LOC.HAIN)
```

**AKD Decision Tree**

```r
names(de) <- c("Sex","Age","AKD","Tarsus","Wing","Mass","LEUC1","HAEM1","PLAS1","SEWA","KENA","VALD","H
ran <- sample(1:nrow(de), 0.9 * nrow(de))
data_train <- de[ran,]
data_test <- de[-ran,]
dtm <- rpart(AKD~., data_train, method="class")
rpart.plot(dtm, compress=TRUE, uniform=TRUE)
```

```
        ┌─────────┐
        │    0    │
        │  0.10   │
        │  100%   │
        └─────────┘
```

57,362,363,364,365,367,368,370,371,372,373,374,375,377,378,379,382,383,384,386,38

```
        ┌─────────┐
        │    1    │
        │  0.58   │
        │   14%   │
        └─────────┘
```

**Tarsus = 45.3,46.2,47.1,48.1,48.2,48.8,48.9,50.8,51.3**

```
  ┌──────┐  ┌──────┐  ┌──────┐
  │  0   │  │  0   │  │  1   │
  │ 0.01 │  │ 0.09 │  │ 1.00 │
  │ 86%  │  │  7%  │  │  8%  │
  └──────┘  └──────┘  └──────┘
```

```r
p <- predict(dtm, data_test, type="class")
confMat <- table(data_test$AKD,p)
accuracy <- sum(diag(confMat))/sum(confMat)
return (accuracy*100)
```

```
## [1] 89.47368
```

**LEUC1 Decision Tree**

```r
names(de) <- c("Sex","Age","AKD","Tarsus","Wing","Mass","LEUC1","HAEM1","PLAS1","SEWA","KENA","VALD","HA
ran <- sample(1:nrow(de), 0.9 * nrow(de))
data_train <- de[ran,]
data_test <- de[-ran,]
dtm <- rpart(LEUC1~., data_train, method="class")
rpart.plot(dtm, compress=TRUE, uniform=TRUE)
```

```r
p <- predict(dtm, data_test, type="class")
confMat <- table(data_test$LEUC1,p)
accuracy <- sum(diag(confMat))/sum(confMat)
return (accuracy*100)
```
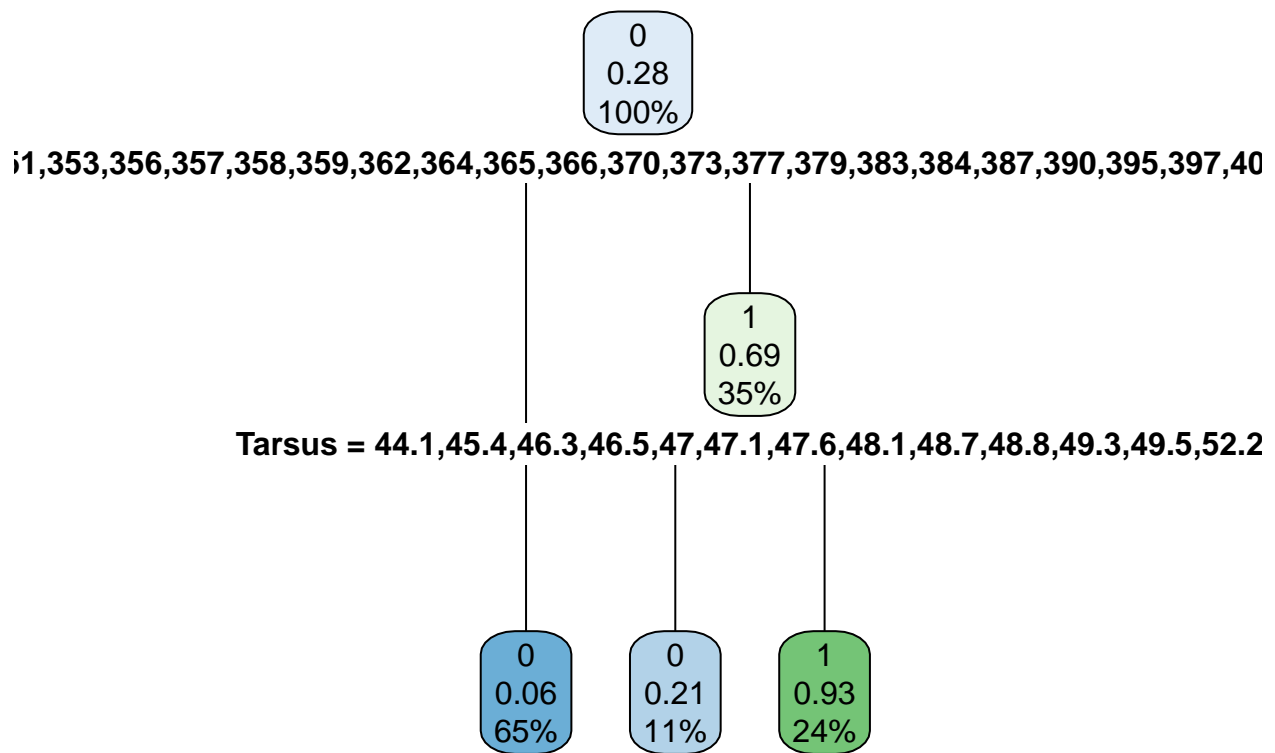
```
## [1] 42.10526
```

**HAEM1 Decision Tree**

```r
names(de) <- c("Sex","Age","AKD","Tarsus","Wing","Mass","LEUC1","HAEM1","PLAS1","SEWA","KENA","VALD","H.
ran <- sample(1:nrow(de), 0.9 * nrow(de))
data_train <- de[ran,]
data_test <- de[-ran,]
dtm <- rpart(HAEM1~., data_train, method="class")
rpart.plot(dtm, compress=TRUE, uniform=TRUE)
```
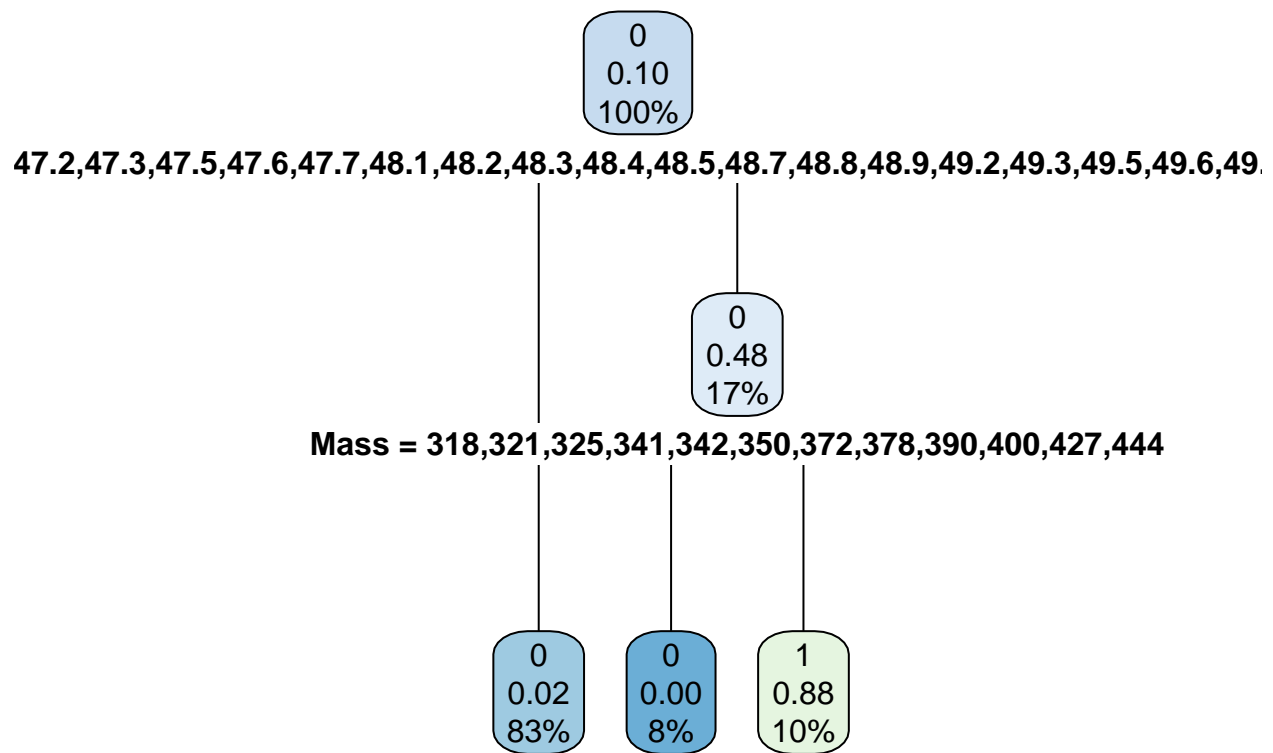
```
p <- predict(dtm, data_test, type="class")
confMat <- table(data_test$HAEM1,p)
accuracy <- sum(diag(confMat))/sum(confMat)
return (accuracy*100)
```

```
## [1] 47.36842
```

**PLAS1 Decision Tree**

```
names(de) <- c("Sex","Age","AKD","Tarsus","Wing","Mass","LEUC1","HAEM1","PLAS1","SEWA","KENA","VALD","HA
ran <- sample(1:nrow(de), 0.9 * nrow(de))
data_train <- de[ran,]
data_test <- de[-ran,]
dtm <- rpart(PLAS1~., data_train, method="class")
rpart.plot(dtm, compress=TRUE, uniform=TRUE)
```

```
        ┌─────────┐
        │    0    │
        │  0.10   │
        │  100%   │
        └─────────┘
```

**47.2,47.3,47.5,47.6,47.7,48.1,48.2,48.3,48.4,48.5,48.7,48.8,48.9,49.2,49.3,49.5,49.6,49.**

```
                        ┌─────────┐
                        │    0    │
                        │  0.48   │
                        │   17%   │
                        └─────────┘
```

**Mass = 318,321,325,341,342,350,372,378,390,400,427,444**

```
   ┌─────────┐  ┌─────────┐  ┌─────────┐
   │    0    │  │    0    │  │    1    │
   │  0.02   │  │  0.00   │  │  0.88   │
   │   83%   │  │    8%   │  │   10%   │
   └─────────┘  └─────────┘  └─────────┘
```

```r
p <- predict(dtm, data_test, type="class")
confMat <- table(data_test$PLAS1,p)
accuracy <- sum(diag(confMat))/sum(confMat)
return (accuracy*100)
```

```
## [1] 84.21053
```